



UNIVERSIDAD
Rafael Belloso Chacín



Revista Electrónica de
Estudios Telemáticos

TÉLÉMATIQUE

RECUPERACIÓN DE INFORMACIÓN: UN ÁREA DE INVESTIGACIÓN EN CRECIMIENTO

Information retrieval: A Growing Area of research

Fernando R. A. Bordignon
Universidad Nacional de Luján, Argentina

Gabriel H. Tolosa
Universidad Nacional de Luján, Argentina

RESUMEN

A partir de la expansión y consolidación de Internet, como medio principal de comunicación electrónica de datos, se ha puesto a disposición de casi toda la humanidad una importante cantidad de información de todo tipo. A los efectos de aprovechar todo este potencial de información, es necesario poseer accesos que permitan que la tarea de recuperación sea efectiva y eficiente en términos de recursos invertidos por los usuarios. Este artículo plantea cuál es el objeto de estudio del área denominada "recuperación de información", en que estado se encuentra y cuales son sus principales líneas de trabajo.

Palabras clave: Recuperación de información, web, motores de búsqueda

ABSTRACT

As a result of the expansion and consolidation of the Internet as the main medium for the transmission of electronic data, a huge amount of information of all kinds has become readily available to humanity. For the purpose of exploiting this information potential it is necessary to have ways of access that would make the information retrieval task both effective and efficient in terms of user resources. This article describes the object of study of "Information Retrieval", its state of the art and main lines of research.

Keywords: Information retrieval, web, search engines

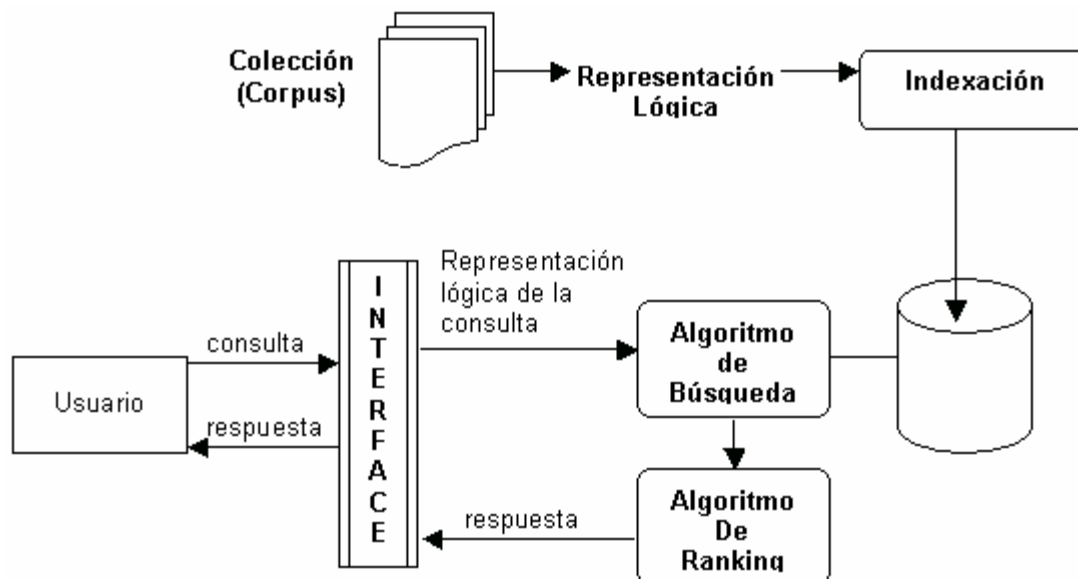


Figura 2 – Arquitectura básica de un SRI

Como podemos observar, se inicia desde un conjunto de documentos de texto, los cuales están compuestos por sucesiones de palabras que forman estructuras gramaticales (por ejemplo, oraciones y párrafos). Tales documentos están escritos en lenguaje natural y expresan ideas de su autor sobre un determinado tema. El conjunto de todos los documentos con los que se trata y sobre los que se deben realizar operaciones de RI se denomina **corpus**, **colección** o **base de datos textual o documental**. Para poder realizar operaciones sobre un corpus, es necesario obtener primero una **representación lógica** de todos sus documentos, la cual puede consistir en un conjunto de términos, frases u otras unidades (sintácticas o semánticas) que permitan – de alguna manera – caracterizarlos. Por ejemplo, la representación de los documentos mediante un conjunto de sus términos se la conoce como “bolsa de palabras” (*bag of words*).

A partir de la representación lógica existe un proceso (**indexación**) que llevará a cabo la construcción de estructuras de datos (normalmente denominadas **índices**) que la almacene y soporte búsquedas eficientes. Es importante destacar que una vez construidos los índices, los documentos del corpus pueden ser eliminados del sistema ya que éste retornará las referencias a los mismos porque cuenta con la información necesaria para hacerlo. En tal caso, el usuario será el encargado de localizar el documento para consultarlo. A los sistemas que funcionan bajo este modelo se los denomina “sistemas referenciales”, en contraste con los que sí almacenan y



representación como imágenes, registro de conversaciones y video. [CLOUGH]

- Desarrollo de Conjuntos (data-sets) de Prueba: A los efectos de evaluar SRI completos o nuevos métodos y técnicas es necesario disponer de juegos de prueba normalizados (corpus con preguntas y respuestas predefinidas, corpus clasificados, etc.). Esta área tiene que ver con la producción tales conjuntos, a partir de diferentes estrategias que permitan reducir la complejidad de la tarea, manejando la dificultad inherente a la carga de subjetividad existente. [GUSTMAN] [SANDERSON]

¿RECUPERACIÓN DE INFORMACIÓN O RECUPERACIÓN DE DATOS?

Muchos usuarios se encuentran familiarizados con el concepto de recuperación de datos (RD), especialmente aquellos que – a menudo – interactúan con sistemas de consulta en bases de datos relacionales ó en registros de alguna naturaleza, como por ejemplo, un registro de los empleados de una organización. Sin embargo, hay diferencias significativas en los conceptos que definen que el tratamiento de las unidades (datos o información) en cada caso sean completamente diferentes.

Básicamente, existen diferencias sustanciales en cuanto a los objetos con que se trata y su representación, la especificación de las consultas y los resultados.

En el área de RD los objetos que se tratan son estructuras de datos conocidas. Su representación se basa en un formato previo definido y con un significado implícito (hay una sintaxis y semántica no ambigua) para cada elemento. Por ejemplo, una tabla en una base de datos que almacena instancias de clientes de una organización posee un conjunto de columnas que definen los atributos de todos los clientes y cada fila corresponde a uno en particular. Nótese que cada elemento (atributo) tiene un dominio conocido y su semántica está claramente establecida. Por otro lado, en el área de RI la unidad u objeto de tratamiento es básicamente un documento de texto – en general – sin estructura.

En cuanto a la especificación de las consultas, en el área de RD se cuenta con una estructura bien definida dada por un lenguaje de consulta que permite su especificación de manera exacta. Las consultas no son ambiguas y consisten en un conjunto de condiciones que deben cumplir los ítems a evaluar para que la misma se satisfaga. Por ejemplo, en el modelo



Como mencionamos, la recuperación de información intenta resolver el problema de encontrar documentos relevantes que satisfagan la necesidad de información del usuario. Sin embargo, se ha planteado la dificultad para llevar a cabo esta tarea debido a la imposibilidad de expresar exactamente tal necesidad. Además, la noción de relevancia es un juicio subjetivo [RIJSBERGEN] y depende de diferentes factores relacionados más cercanamente con el usuario. La relevancia de un documento respecto a un *query* se refiere a cuánto el primero responde al segundo. De igual manera, luego el usuario evalúa qué tanto, es decir, en qué medida, se satisface su necesidad de información [KORFHAGE].

Es por ello, que se plantea la relevancia como similitud, para poder comparar documentos con consultas y – bajo ciertos criterios – definir una medida de distancia entre ambos. Por lo tanto, se puede plantear la idea que “un documento es relevante a una consulta si son similares”, donde la medida de similitud puede estar basada en diferentes criterios (coincidencias de términos, significado de éstos, frecuencia de aparición de términos y distribución del vocabulario, entre otros).

Martínez Méndez y otros [MARTINEZ] resaltan la dificultad para determinar la relevancia o no de un documento respecto de una consulta. Plantean – por ejemplo – que dos personas pueden juzgar un mismo documento de diferente manera y que es difícil establecer los criterios para la evaluación de la relevancia. Finalmente, mencionan la idea de relevancia parcial, es decir, cuando solo una parte del documento se considera relevante.

Por otro lado, como el *query* no describe exactamente la necesidad de información del usuario, algunos autores [KORFHAGE] definen el concepto de “pertinencia”, donde se incluyen las restricciones impuestas por el SRI. Este concepto está relacionado con la utilidad del documento para el usuario [MARTINEZ], de acuerdo a la necesidad de información original que guió su búsqueda, independientemente si es en parte o todo el documento.

Sin embargo – y a pesar de las dificultades para determinarla – el concepto genérico de relevancia es aceptado ampliamente por la comunidad de RI para evaluar la respuesta de un SRI respecto de una consulta de un usuario, la cual – como ya mencionamos – surge a partir de una necesidad de información.



Si bien es el primer modelo desarrollado y aún se lo utiliza, no es el preferido por los ingenieros de software para sus desarrollos. Existen diversos puntos en contra que hacen que cada día se lo utilice menos y – además – se han desarrollado algunas extensiones, bajo el nombre modelo booleano extendido [WALLER] [SALTON_b], que tratan de mejorar algunos puntos débiles.

b) MODELO VECTORIAL

Este modelo fue planteado y desarrollado por Gerard Salton [SALTON_c] y – originalmente – se implementó en un SRI llamado SMART. Aunque el modelo posee más de treinta años, actualmente se sigue utilizando debido a su buena performance en la recuperación de documentos.

Conceptualmente, este modelo utiliza una matriz documento–término que contiene el vocabulario de la colección de referencia y los documentos existentes. En la intersección de un término t y un documento d se almacena un valor numérico de importancia del término t en el documento d ; tal valor representa su *poder de discriminación*. Así, cada documento puede ser visto como un vector que pertenece a un espacio n -dimensional, donde n es la cantidad de términos que componen el vocabulario de la colección. En teoría, los documentos que contengan términos similares estarán a muy poca distancia entre sí sobre tal espacio. De igual forma se trata a la consulta, es un documento más y se la mapea sobre el espacio de documentos. Luego, a partir de una consulta dada es posible devolver una lista de documentos ordenados por distancia (los más relevantes primero). Para calcular la semejanza entre el vector consulta y los vectores que representan los documentos se utilizan diferentes fórmulas de distancia, siendo la más común la del coseno.

Obsérvese el siguiente ejemplo donde se representa a un documento d y a una consulta c :

Documento: “*La República Argentina ha sido nominada para la realización del X Congreso Americano de Epidemiología en Zonas de Desastre. El encuentro se realizará...*”

Consulta: “*argentina congreso epidemiología*”

	argentina	...	congreso	epidemiología	...
--	-----------	-----	----------	---------------	-----

elementos multimedia, mayor precisión sobre el ámbito de la consulta y demás).

La estructura de los documentos a indexar está dada por marcas o etiquetas, siendo los estándares más utilizados el SGML (*Standard General Markup Language*), el HTML (*HyperText Markup Language*), el XML (*eXtensible Markup Language*) y LATEX.

Al poseer la descripción de parte de la estructura de un documento es posible generar un grafo sobre el que se navegue y se respondan consultas de distinto tipo, por ejemplo:

- Por estructura: *¿Cuáles son las secciones del segundo capítulo?*
- Por metadatos o campos: *Documentos de “Editorial UNLu” editados en 1998*
- Por contenido: *Término “agua” en títulos de secciones*
- Por elementos multimedia: *Imágenes cercanas a párrafos que contengan “Bosch”*

Para Baeza-Yates existen dos modelos en esta categoría “nodos proximales” [NAVARRO] y “listas no superpuestas” [BURKOWSKI]. Ambos modelos se basan en almacenar las ocurrencias de los términos a indexar en estructuras de datos diferentes, según aparezcan en algún elemento de estructura (región) o en otro como capítulos, secciones, subsecciones y demás. En general, las regiones de una misma estructura de datos no poseen superposición, pero regiones en diferentes estructuras sí se pueden superponer. Los tipos de consultas soportados son simples:

- Seleccione una región que contenga una palabra dada
- Seleccione una región X que no contenga una región Y
- Seleccione una región contenida en otra región

Sobre una estructura tipo libro un ejemplo de consulta sería:

[subsección[+] CONTIENE “tambo”]

Como respuesta el SRI buscaría subsecciones y sub-subsecciones que contengan el término “tambo”.

Cabe mencionar que algunos motores de búsqueda de Internet ya utilizan ciertos elementos de la estructura de un documento – por ejemplo, los



- [SALTON] Salton, G. (1971) (editor). **The SMART Retrieval System – Experiments in Automatic Document Processing**. Ed. Prentice Hall Inc. Englewood Cliffs, NJ.
- [SMEDT] Smedt, K. Liseth, A. Hassel, M. y Dalianis, H. (2005). **How short is good? An evaluation of automatic summarization**. En Holmboe, H. (ed.) **Nordisk Sprogteknologi 2004**. Årbog for Nordisk Språkteknologisk Forskningsprogram 2000-2004, pp 267-287.
- [WADE] Wade, C. y Allan, J., (2005) **Passage Retrieval and Evaluation**, CIIR Technical Report.
- [WALLER] Waller, W. G. y Kraft, D. H. (1979) **A mathematical model for a weighted Boolean retrieval system**". Information Processing and Management, 15(5):235-245. 1979.